# Chemical Space is Infinite: How can one scale to infinity while still being usable/useful?
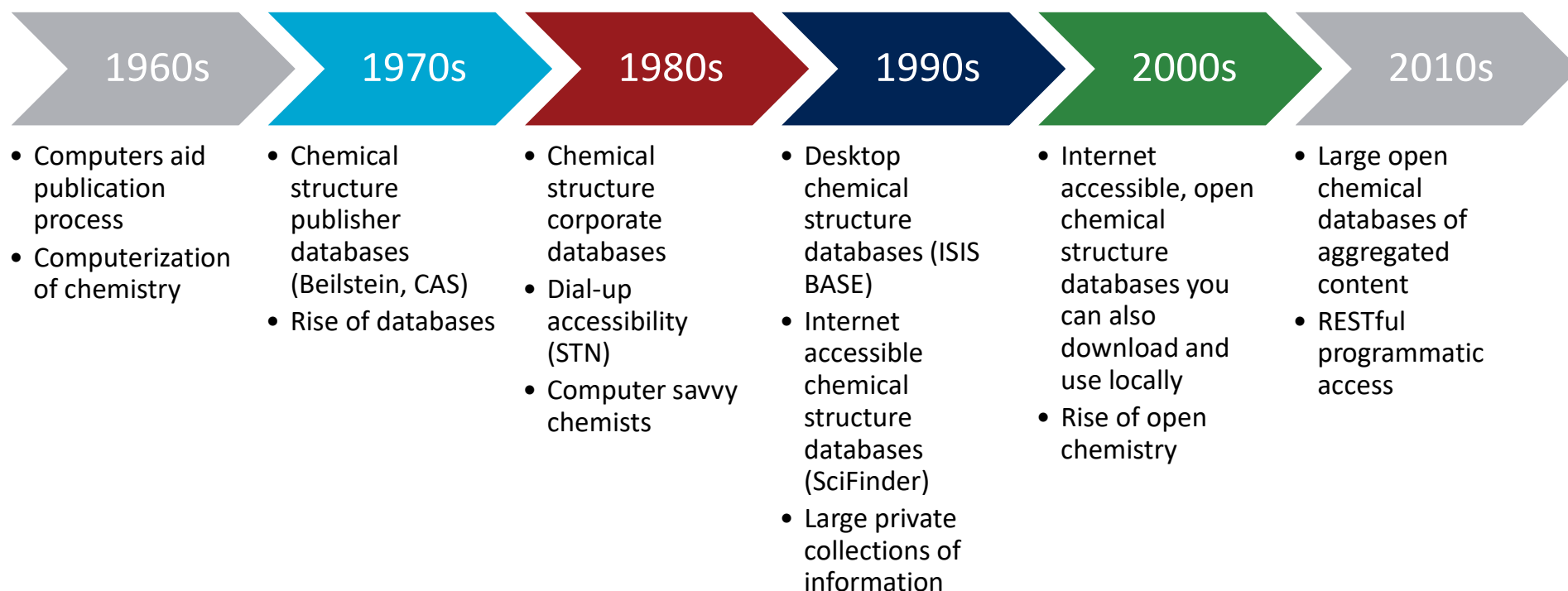
Evan Bolton, Ph.D.  -  Program Head of Chemistry

**NIH** U.S. National Library of Medicine
National Center for Biotechnology Information

# An evolution of chemical structure databases

| 1960s | 1970s | 1980s | 1990s | 2000s | 2010s |
|-------|-------|-------|-------|-------|-------|

- **1960s**
  - Computers aid publication process
  - Computerization of chemistry

- **1970s**
  - Chemical structure publisher databases (Beilstein, CAS)
  - Rise of databases

- **1980s**
  - Chemical structure corporate databases
  - Dial-up accessibility (STN)
  - Computer savvy chemists

- **1990s**
  - Desktop chemical structure databases (ISIS BASE)
  - Internet accessible chemical structure databases (SciFinder)
  - Large private collections of information

- **2000s**
  - Internet accessible, open chemical structure databases you can also download and use locally
  - Rise of open chemistry

- **2010s**
  - Large open chemical databases of aggregated content
  - RESTful programmatic access

NIH — U.S. National Library of Medicine
National Center for Biotechnology Information

NCBI

*Imagine you wanted to make a modern scientific resource, what do you need to focus on?*

Use of persistent research identifiers

Data use cases explicitly considered

Standards-based approaches

Explicit data licensing (e.g., CC-BY 4.0, CC0)

# 2020s
# Cloud-first,
# Mobile-first,
# Machine-first,
# FAIR-first, Open-first

Receive cloud-based data

Make data accessible within cloud

UI/UX from a device screen size agnostic perspective

Use of controlled vocabulary and machine interpretable statements

Sufficient meta data to reproduce the science

**FAIR** means "**F**ully **AI R**eady" .. also means "**F**indable", "**A**ccessible", "**I**nteroperable", "**R**eusable"

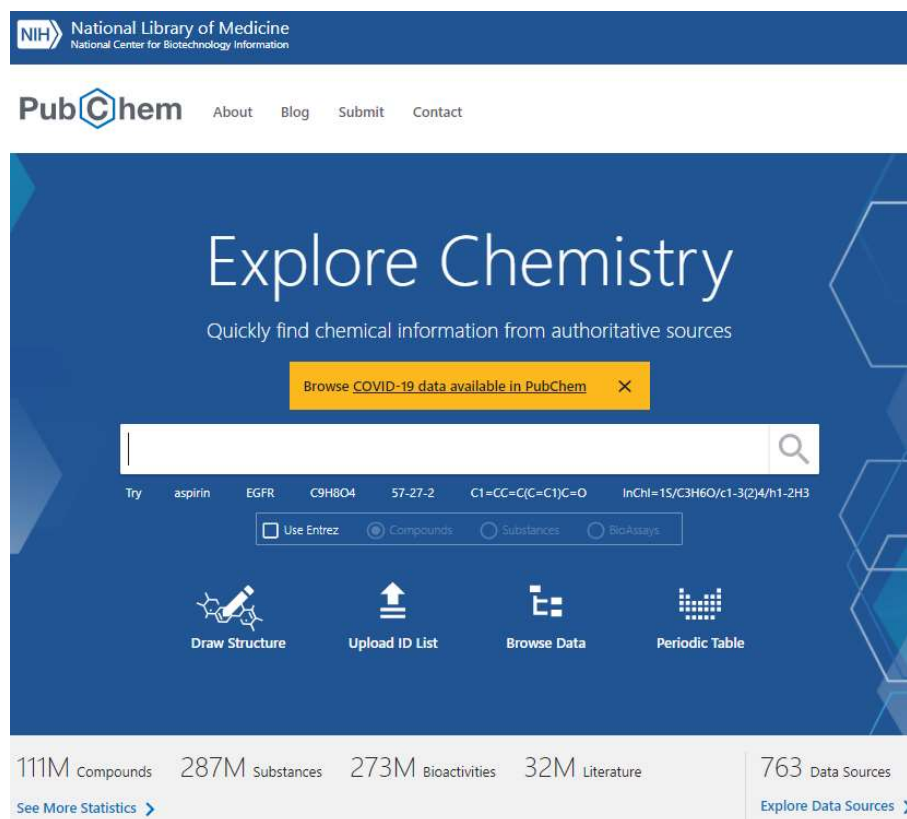Are we ready for ultra-large chemical databases?

- InChIKey, will it still work?
- Users, will they know how to use?
- Is interactivity a thing of the past?
- Are the possibilities so great that it all just seems random?

Image credit:
https://ilekh.com/wp-content/uploads/2014/08/time-for-change.jpg

# PubChem is a data repository

- World's largest collection of freely accessible chemical information.

- Helps researchers make sense of the biological roles and health effects of chemicals on human health and the environment.
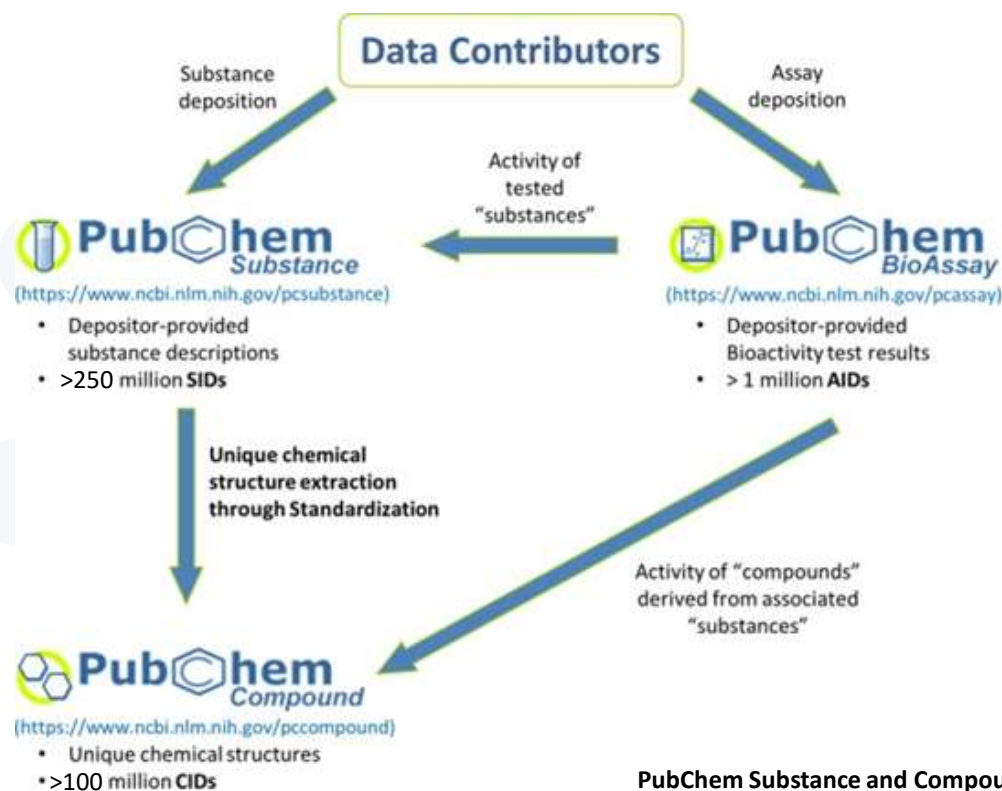
U.S. National Library of Medicine
National Center for Biotechnology Information

Chemical substances and bioactivities .. with select annotation



https://pubchem.ncbi.nlm.nih.gov/

# Two primary archival databases

Compound is derived from Substance



**Data Contributors**

Substance deposition

Assay deposition

Activity of tested "substances"

**PubChem** *Substance*
(https://www.ncbi.nlm.nih.gov/pcsubstance)
- Depositor-provided substance descriptions
- >250 million **SIDs**

**PubChem** *BioAssay*
(https://www.ncbi.nlm.nih.gov/pcassay)
- Depositor-provided Bioactivity test results
- > 1 million **AIDs**

Unique chemical structure extraction through Standardization

Activity of "compounds" derived from associated "substances"

**PubChem** *Compound*
(https://www.ncbi.nlm.nih.gov/pccompound)
- Unique chemical structures
- >100 million **CIDs**

NIH⟩ U.S. National Library of Medicine
National Center for Biotechnology Information

**PubChem** *Substance*

↓

- ❖ **Validate chemical contents**
  - Atoms defined/real
  - Implicit hydrogen
  - Functional group
  - Atom valence

⇩

- ❖ **Normalize representations**
  - Tautomer invariance
  - Aromaticity detection
  - Stereochemistry
  - Explicit hydrogen

⇩

- ❖ **Calculate**
  - 2-D depiction coordinates
  - Molecular properties
  - Chemical descriptors

⇩

- ❖ **Additional processing for mixtures**
  - Isolate covalent units
  - Neutralize (by ± H⁺ or e⁻)
  - Reprocess
  - Detect unique components

⇩

**PubChem** *Compound*

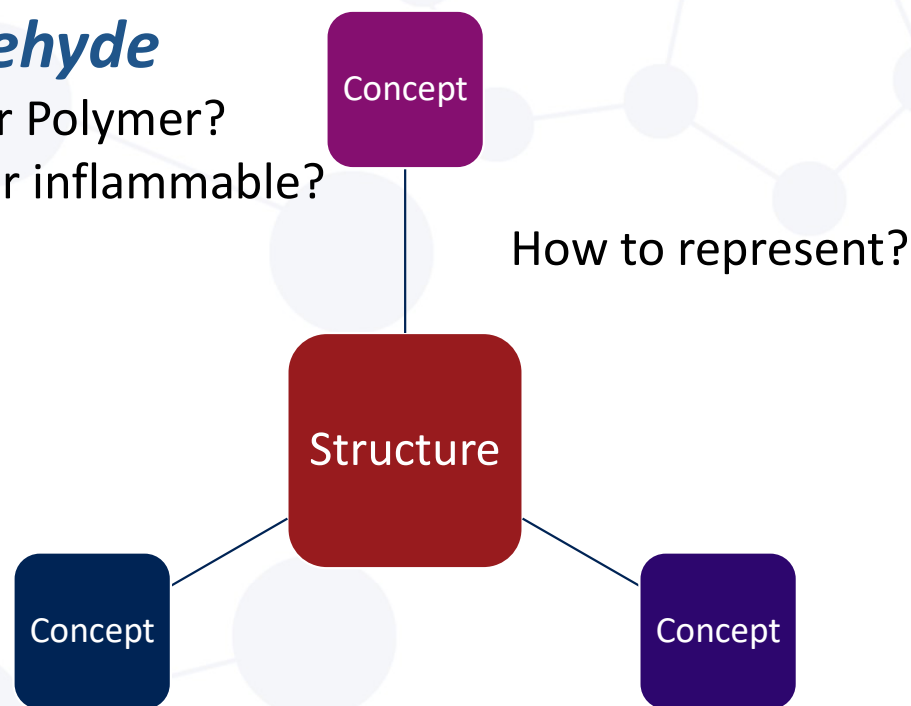Many to many links abound in PubChem within a given collection

# Many to many relationships

*Carbon*

Element?
Coal?
Diamond?
Methane?

*Formaldehyde*

Gas or Liquid or Polymer?
Liquid: flammable or inflammable?

How to represent?

*Gleevec*

Salt?
Hydrate?
Free base?

Structure — Concept — Structure — Structure — Structure

Concept — Structure — Concept — Concept

# Integrating many resources

- Over 700 contributing resources
  - Unified search engine
    - Flexible handling of query types
    - Includes chemical input extensions and sketcher
  - Many collection types
    - Compounds / Substances
    - Proteins / Genes / *Pathways*
    - BioAssays
    - Literature / Patents

## 15.2 Springer Nature References

16,041 items  View More Rows & Details  ⤢                    ⬇ Download

| | | SORT BY ⇕ Relevance | | ⌄ |
|---|---|---|---|---|

| Thumbnail | Title | Publication Name | Publication Date | PMID |
|---|---|---|---|---|
| Journal of Materials Science | Magnetism and white-light-emission bifunctionality simultaneously assembled into flexible Janus nanofiber via electrospinning | Journal of Materials Science | 2015 | |
| Materials in Electronics | Facile electrospinning construction and characteristics of coaxial nanobelts with simultaneously tunable magnetism and color-tuned photoluminescence bifunctionality | Journal of Materials Science: Materials in Electronics | 2015 | |
| JOURNAL OF NANOPARTICLE RESEARCH | One-pot facile electrospinning construct of flexible Janus nanofibers with tunable and enhanced magnetism–photoluminescence bifunctionality | Journal of Nanoparticle Research | 2015 | |
| Materials in Electronics | A novel scheme to obtain tunable fluorescent colors based on electrospun composite nanofibers | Journal of Materials Science: Materials in Electronics | 2014 | |

# Integrating publisher provided metadata

- Major publishers provide chemical-DOI associations
  - Thieme
  - Springer Nature
- Publishers provide document level metadata
  - CrossRef, PubMed, SciGraph, (Agricola)

15.7 Chemical Co-Oc...
15.8 Chemical-Disease
15.9 Chemical-Gene Co-Occurrences in Literature

Showing 3 of 25 View More Co-Occ...
Showing 3 of 25 View More Co-Occu...
Showing 3 of 100 View More Co-Occurrence and Evidence Data
Download

| Chemical | Evidence from |
|---|---|
| Salicylic Acid CID 338 | 359 articles<br>Comparative ph... acid on photosy...<br>PMID 29751250; PI<br>Name matches: sa...<br>Terahertz (6-15 ... Vibrations in Be...<br>PMID 25909770; A<br>Name matches: 2-...<br>Cloning and cha... methyltransfera... 'Yelloween').<br>PMID 26600510; G...<br>Name matches: sa... |
| Benzaldehyde CID 240 | 188 articles<br>A new enzymati... cinnamic acid to...<br>PMID 30417393; P...<br>Name matches: be...<br>Benzaldehyde in...<br>PMID 27041300; F...<br>Name matches: be...<br>Simultaneous q... pharmaceuticals... followed by liqu...<br>PMID 27495371; Jo...<br>Name matches: be... |

| Disease | Evidence fro... |
|---|---|
| Drug-Related Side Effects And Adverse Reactions | 207 articles<br>Identificatio... highly poter...<br>PMID 3064269...<br>Name matches...<br>Bulk Organo... Lactide and ...<br>PMID 3096427...<br>Name matche...<br>Rapid gener... potent, selec... prodrug stu...<br>PMID 3133630...<br>Name matche... |
| Neoplasms | 115 articles<br>Comparative... derivatives a... inflammator...<br>PMID 2592762...<br>Name matches...<br>Hydroxamic... Glioma and ...<br>PMID 2608812...<br>Name matches...<br>Naturally oc... inhibiting hi...<br>PMID 2850619... |

| Gene | Evidence from All Time ∨ |
|---|---|
| Tyrosinase | 33 articles   ⬇ Download CSV  View in PubMed ⎘<br>Tyrosinase inhibitory effect of benzoic acid derivatives and their structure-activity relationships.<br>PMID 20476840; Journal of enzyme inhibition and medicinal chemistry 2010 Dec; 25(6):812-817<br>Name matches: **tyrosinase** benzoic acid<br>An optical test strip for the detection of benzoic acid in food.<br>PMID 22164018; Sensors (Basel, Switzerland) 2011 ; 11(8):7302-7313<br>Name matches: **tyrosinase** benzoic acid<br>Tyrosinase biosensor for benzoic acid inhibition-based determination with the use of a flow-batch monosegmented sequential injection system.<br>PMID 22817942; Talanta 2012 Jul; 96(?):147-152<br>Name matches: **tyrosinase** benzoic acid |
| Monocarboxylic Acid Transporter | 18 articles   ⬇ Download CSV  View in PubMed ⎘<br>Absorption of benzoic acid in segmental regions of the vascularly perfused rat small intestine preparation.<br>PMID 11717172; Drug metabolism and disposition: the biological fate of chemicals 2001 Dec; 29(12):1539-1547<br>Name matches: **monocarboxylic acid transporter** benzoic acid<br>Transepithelial transport of artepillin C in intestinal Caco-2 cell monolayers.<br>PMID 16004960; Biochimica et biophysica acta 2005 Jul; 1713(2):138-144<br>Name matches: **monocarboxylic acid transporter** benzoic acid<br>Uptake of 4-chloro-2-methylphenoxyacetic acid (MCPA) from the apical membrane of Caco-2 cells by the monocarboxylic acid transporter.<br>PMID 18096194; Toxicology and applied pharmacology 2008 Mar; 227(3):325-330<br>Name matches: **monocarboxylic acid transporter** benzoic acid |

▸ from PubChem          ▸ from PubChem

- Using PubMed corpus
- Mine text title/abstract
- Find all chemical, disease, gene/protein mentions
- Compute histogram and provide top-N
- Evidence clearly stated
  - name, PMIDs, ..
  - Downloadable
- Handles all nine combinations of Chemical, Gene/Protein, Disease

Many precomputed relationships and analyses exist to further integration and interpretation by users and machines

PubChem Co-occurrence
*'Knowledge Panel'* displays

# The chemical information ~~understanding~~ divide

**Human understanding**

Depictions
Schemes
Table
Text

**Human Intent**

**Computer understanding**

Explicit
Complete
Annotated
Interpreted

# PubChem RDF-based Linked Data

- Describes relationships between PubChem data

- Machine-readable information triples

- Organized like layers of an onion (just take what you need)

- Uses ontologies and vocabulary description

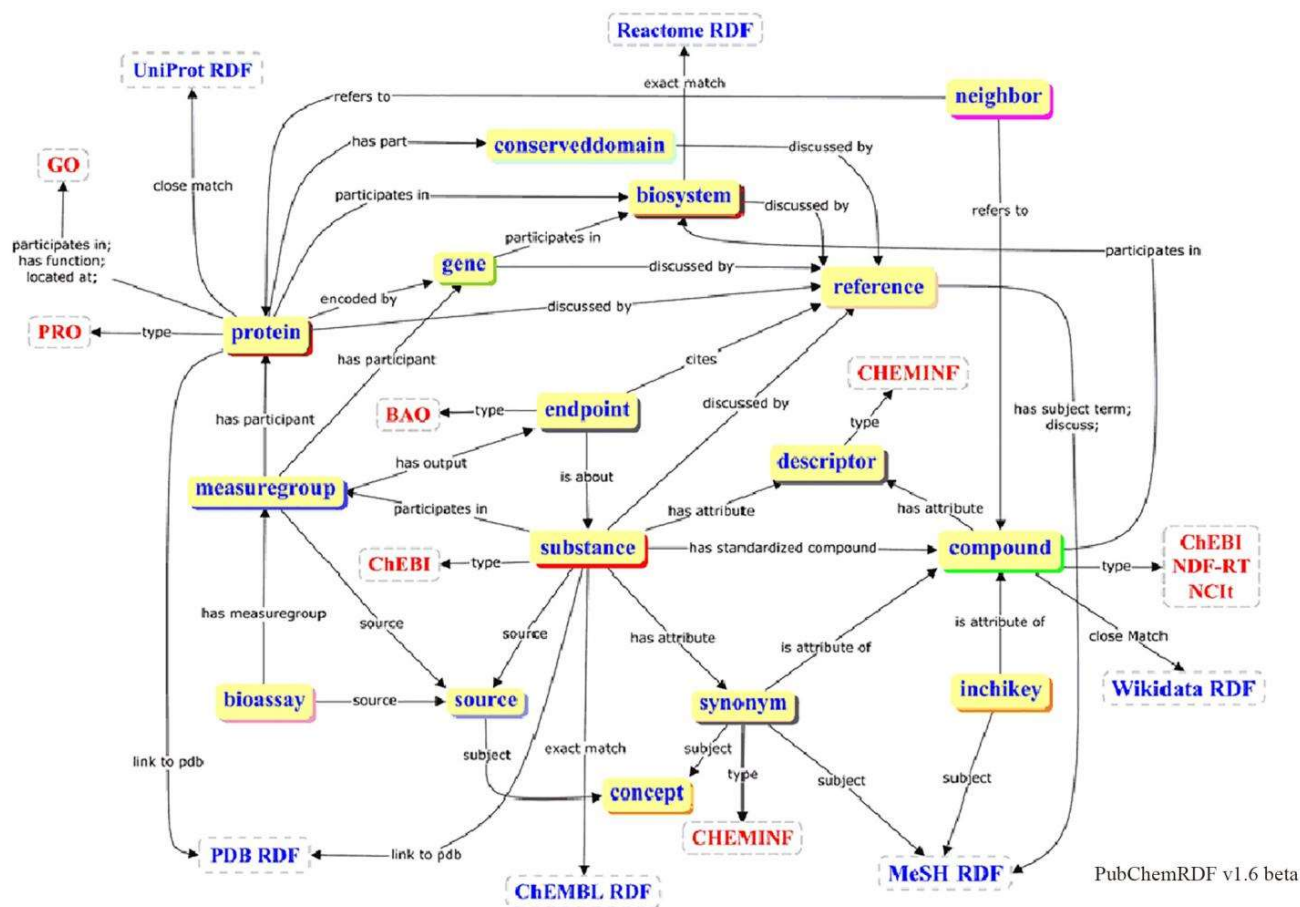- Only a core set of PubChem content
  - Expanding coverage



Figure 1. Color-coded diagram showing a high-level overview of the PubChemRDF semantic relationships.

# PubChem RDF-based Linked Data

- Describes relationships between PubChem data
- Machine-readable information triples
- Organized like layers of an onion (just take what you need)
- Uses ontologies and vocabulary description
- Only a core set of PubChem content
- Working to add more

NIH | U.S. National Library of Medicine
National Center for Biotechnology Information

# PubChemRDF Statistics

**Total number of triples: 73,443,150,086**

Last updated on 01-31-2020

| Prefix/Namespace | Total number of triples | Total number of subjects |
|---|---|---|
| compound https://rdf.ncbi.nlm.nih.gov/pubchem/compound/ | Non-neighboring links: 2,466,218,961 2D neighboring links: 29,325,920,096 3D neighboring links: 32,373,792,809 | 102,429,168 |
| substance https://rdf.ncbi.nlm.nih.gov/pubchem/substance/ | 1,775,934,055 | 388,300,240 |
| descr https://rdf.ncbi.nlm.nih.gov/pubchem/descriptor/ | 5,624,993,164 | 2,607,822,222 |
| inchikey https://rdf.ncbi.nlm.nih.gov/pubchem/inchikey/ | 306,689,898 | 102,127,699 |
| syno https://rdf.ncbi.nlm.nih.gov/pubchem/synonym/ | 480,509,253 | 181,804,989 |
| bioassay https://rdf.ncbi.nlm.nih.gov/pubchem/bioassay/ | 98,335 | 23400 |
| measuregroup https://rdf.ncbi.nlm.nih.gov/pubchem/measuregroup/ | 248,048,134 | 1,090,853 |

… and more …

NCBI

# PubChem FTP Site

- Different file formats
- 7 TBs of data
- Compound, Substance, BioAssay, Bioactivities
- Lots of extras
  - Special data sets
  - Link files to patents, PubMed, …
  - Target
  - RDF
  - Specifications

https://ftp.ncbi.nlm.nih.gov/pubchem/

## Index of /pubchem

Everything is downloadable in bulk

| Name | Last modified | Size |
|---|---|---|
| Parent Directory | | - |
| Bioassay/ | 2019-04-17 11:31 | - |
| Compound/ | 2019-07-12 02:59 | - |
| Compound_3D/ | 2019-07-02 08:04 | - |
| Other/ | 2019-04-17 12:40 | - |
| RDF/ | 2019-09-23 19:46 | - |
| Substance/ | 2019-04-17 12:24 | - |
| Target/ | 2017-03-09 19:48 | - |
| data_spec/ | 2019-04-17 12:53 | - |
| presentations/ | 2016-03-04 12:31 | - |
| publications/ | 2019-04-17 12:43 | - |
| specifications/ | 2019-04-17 12:53 | - |
| README | 2016-11-04 11:28 | 1.7K |

ftp://ftp.ncbi.nlm.nih.gov/pubchem/
https://ftp.ncbi.nlm.nih.gov/pubchem/

==Annotation-based data access==

Database | Open Access | Published: 09 August 2019

## PUG-View: programmatic access to chemical annotations integrated in PubChem

Sunghwan Kim, Paul A. Thiessen, Tiejun Cheng, Jian Zhang, Asta Gindulyte & Evan E. Bolton ✉

*Journal of Cheminformatics* 11, Article number: 56 (2019) | Download Citation ⬇

557 Accesses | 15 Altmetric | Metrics »

PMID: 31399858
PMCID: PMC6324075
DOI: 10.1186/s13321-019-0375-2

==Archive-based data access==

## An update on PUG-REST: RESTful interface for programmatic access to PubChem FREE

Sunghwan Kim, Paul A Thiessen, Tiejun Cheng, Bo Yu, Evan E Bolton ✉ Author Notes

*Nucleic Acids Research*, Volume 46, Issue W1, 2 July 2018, Pages W563–W570,
https://doi.org/10.1093/nar/gky294

Published: 30 April 2018 Article history ▾

PMID: 29718389
PMCID: PMC6030920
DOI: 10.1093/nar/gky294

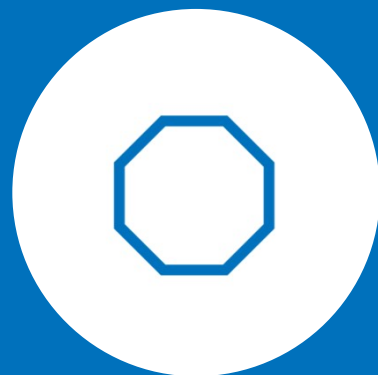# Content is programmatically accessible

==Everything is analyzable or downloadable in pieces interactively and programmatically .. capacity is needed at scale==

# What would happen if all these counts increased 10x?

# PubChem Data Counts

| Data Collection | Live Count | Description |
|---|---|---|
| Compounds | 111,458,063 | Unique chemical structures extracted |
| Substances | 287,046,030 | Information about chemical entities p |
| BioAssays | 1,229,043 | Biological experiments provided by Pu |
| Bioactivities | 273,300,136 | Biological activity data points reporte |
| Genes | 91,340 | Gene targets tested in PubChem BioA |
| Proteins | 99,361 | Protein targets tested in PubChem Bio |
| Taxonomy | | Organisms of targets tested in PubCh |
| Pathways | 237,775 | Interactions between chemicals, gene |
| Literature | 31,753,737 | Scientific publications with links in Pu |
| Patents | 24,824,605 | Patents with links in PubChem |
| Data Sources | 762 | Organizations contributing data to Pu |

Short answer .. PubChem would implode .. it does not scale to such a level .. a (complete) rethink would be needed

# Contemplating chemical infinity

# What holds us back from chemical infinity?

One or more of the following:

1. Hardware
2. Software
3. Money
4. Time
5. Use cases

# Very basic operations of a chemical structure database

- Query by structure
  - **Identity** – scales as N
  - **Similarity** – scales as N
  - **Substructure** – scales as N-ish
    (depends on the algorithm and tradeoffs)
- **Sort** results – scales as N * ln N
- **Filter** results – scales as N
- **Retrieve** results – scales as N

(assumes a full scan needed, many optimizations can be applied, completely ignores analysis beyond filtering)

Considering only

| Database Size | N * ln N |
|---|---|
| 1M | 10M |
| 10M | 100M |
| 100M | 1B |
| 1B | 10B |
| 10B | 100B |
| 100B | 1T |
| 1T | 10T |
| 10T | 100T |

M=million, B=billion, T=trillion

NCBI

==Can we implement more memory efficient approaches?==

Consider the economics of a feature-less chemical structure database in the cloud

- Assuming
  - 100 bytes per structure
  - vCPU w/ 2GB memory @ $0.025/hour
  - Storage 1GB @ $0.08/month
  - Minimal I/O used (else $$)
  - All in memory (for speed!)
  - Capacity for only one query at a time

| Database Size | Storage | vCPU needed | Yearly cost |
|---|---|---|---|
| 1M | 100MB | 1 | $ 219 |
| 10M | 1GB | 1 | $ 220 |
| 100M | 10GB | 5 | $ 1,105 |
| 1B | 100GB | 50 | $ 11,046 |
| 10B | 1PB | 500 | $ 110,460 |
| 100B | 10PB | 5,000 | $ 1,104,600 |
| 1T | 100PB | 50,000 | $ 11,046,000 |
| 10T | 1EB | 500,000 | $ 110,460,000 |

(many tradeoffs can be applied to change the economics, ignores many aspects needed to run at scale, a full featured chemical search system likely 10-100 times more expensive before optimization, analysis costs extra)
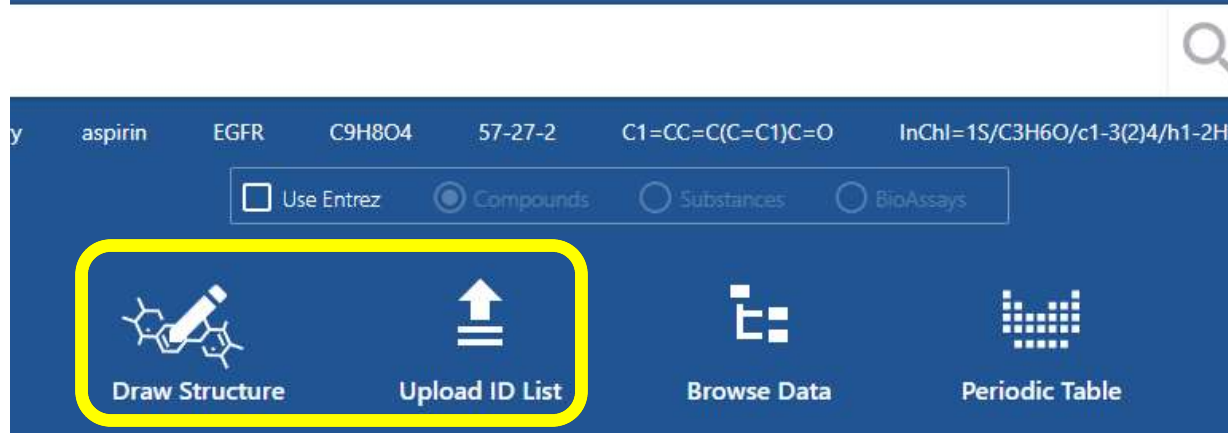
# Practical considerations

(beyond cost)

- How does the workflow change for users?
  - What can a user do with a large selection of results?
  - Does one save a list of 1.5B structures to come back later and analyze more?
  - Will speed be sufficient when store vs. compute-on-the-fly becomes a serious consideration?
  - What decision-making analysis will be useful to users?

Image credit:
https://sassofia.com/wp-content/uploads/2019/03/Untitled-1.jpg

# Explore Chemistry

Quickly find chemical information from authoritative sources

aspirin    EGFR    C9H8O4    57-27-2    C1=CC=C(C=C1)C=O    InChI=1S/C3H6O/c1-3(2)4/h1-2H3

☐ Use Entrez    ⊙ Compounds    ○ Substances    ○ BioAssays

**Draw Structure**    **Upload ID List**    **Browse Data**    **Periodic Table**

How do you import data?
How much can you import?
How do you use large inputs on database side?

# PubChem Search

https://pubchem.ncbi.nlm.nih.gov/

- Single box, many query types
  - Chemical name, CAS#
  - Gene symbol/name
  - Molecular Formula
  - SMILES, InChI
  - SMARTS (substructure)
  - ...

- Draw a structure
  - Or upload a file

- Chemical Search by
  - Identity, similarity, substructure, superstructure, mol. formula

- Upload an ID list

- Many collections
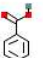  - Compounds, Substances, BioAssays, Genes, Proteins, Pathways, Literature, Patents

# PubChem Search

https://pubchem.ncbi.nlm.nih.gov/

- Single box, many query types
  - Chemical name, CAS#
  - Gene symbol/name
  - Molecular Formula
  - SMILES, InChI
  - SMARTS (substructure)
  - ...

- Draw a structure
  - Or upload a file

- Chemical Search by
  - Identity, similarity, substructure, superstructure, mol. formula

- Upload an ID list

- Many collections
  - Compounds, Substances, BioAssays, Genes, Proteins, Pathways, Literature, Patents

Compounds (771,377) | Substances (453,445) | Genes (1) | Proteins (29) | Pathways (5) | BioAssays (5,230) | Literature (9,886)

Patents (2,296)

Searching chemical names and synonyms including IUPAC names and InChIKeys across the compound collection. Note that annotations text from compound summary pages is not searched. Read More...

771,377 results    Filters

**Choose Sort Options** ✕

- Relevance
- Annotation Record Count
- Compound CID
- Complexity
- H-Bond Donor Count
- H-Bond Acceptor Count
- Heavy Atom Count
- Molecular Weight
- Polar Area
- Rotatable Bond Count
- XLogP
- Create Date

Benzoic Acid; 65-85-0; D... Carboxybenzene; ...
Compound CID: 243
MF: $C_7H_6O_2$   MW: 122.12g/mol
InChIKey: WPYMKLBDIGXBTP-UH...
IUPAC Name: benzoic acid
Create Date: 2004-09-16

Summary   Similar Structures Search   Related

4-(Quinolin-2-Ylmethoxy... Quinolinylmethoxy)Benz... Ylmethyloxy)Benzoic Acid... Acid; ...
Compound CID: 9838553
MF: $C_{17}H_{13}NO_3$   MW: 279.29g/mol
InChIKey: N...
IUPAC Nam...
Create Date...

Summary   Similar Structur...

Are you allowed to filter?
How do you do it?
How will it affect performance?
How will it scale?

4-[(Nitrooxy)Methyl]Benzoic Acid; 258278-55-6; 4-(Nitrooxymethyl)Benzoic Acid; 4-Nitrooxymethylbenzoic Acid; SCHEMBL186313; ...
Compound CID: 9855737

Download CSV

Search in Entrez

ACTIONS ON RESULTS WITH ID TYPE:
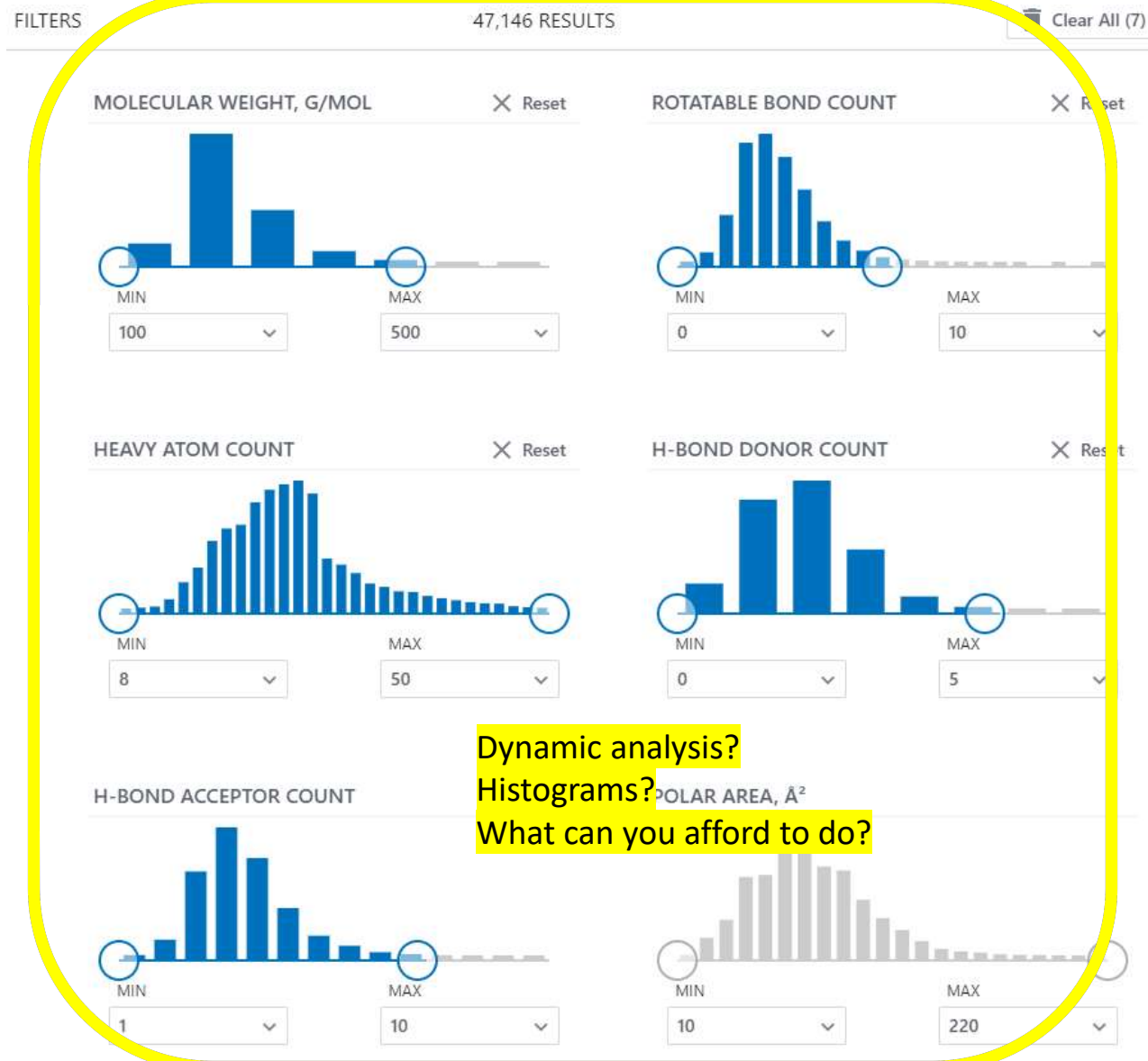CID - Compounds

Push to Entrez
Save for Later
Linked Data Sets

- **Single text box, many query types**
  - Chemical name, CAS#
  - Gene symbol/name
  - Molecular Formula
  - SMILES, InChI
  - SMARTS (substructure)
  - ...
- **Draw a structure**
  - Or upload a file
- **Chemical Search by identity, similarity, substructure, superstructure, mol. formula**
- **Upload an ID list**
- **Many collections**
  - Compounds, Substances, BioAssays, Genes, Proteins, Pathways, Literature, Patents

- **Single text box, many query types**
  - Chemical name, CAS#
  - Gene symbol/name
  - Molecular Formula
  - SMILES, InChI
  - SMARTS (substructure)
  - …

- **Draw a structure**
  - Or upload a file

- **Chemical Search by identity, similarity, substructure, superstructure, mol. formula**

- **Upload an ID list**

- **Many collections**
  - Compounds, Substances, BioAssays, Genes, Proteins, Pathways, Literature, Patents

# PubChem3D

## Research Article

### Similar compounds versus similar conformers: complementarity between PubChem 2-D and 3-D neighboring sets
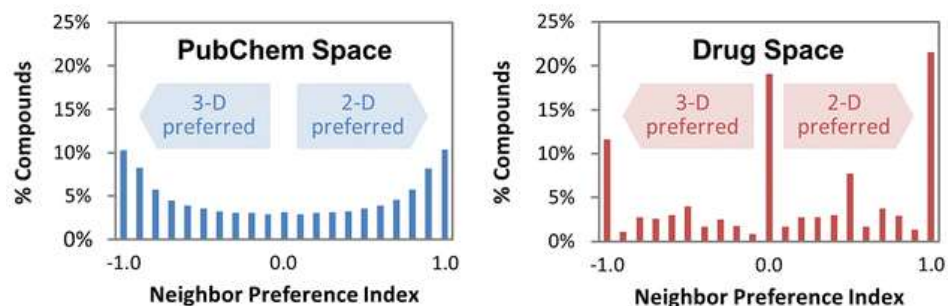
PubChem is a public repository for biological activities of small molecules. For the efficient use of its vast amount of chemical information, PubChem performs 2-dimensional (2-D) and 3-dimensional (3-D) neigh...

Sunghwan Kim, Evan E. Bolton and Stephen H. Bryant

*Journal of Cheminformatics* 2016 8:62
Published on: 4 November 2016

> Full Text    > PDF



## Research Article

### PubChem structure–activity relationship (SAR) clusters

Developing structure–activity relationships (SARs) of molecules is an important approach in facilitating hit exploration in the early stage of drug discovery. Although information on millions of compounds and

Thematic Series of ten articles
https://www.biomedcentral.com/collections/pubchem3d

---

**What about chemical similarity?**
**Will it have meaning in an ultra-large DB?**
**Can you afford to compute (and store?) a fingerprint?**
**Dare you try 3-D similarity? (can you afford the GPUs?)**

---

- Computationally generated 3-D structure
- Set of ten diverse conformers
- 3-D display widget
- 3-D similarity measure
  - Similar shape and protein binding features
  - Compliments 2-D similarity
- Downloadable data
  - SDF, precomputed similarity

# Classification Browser

**Browse Data**

- Multiple hierarchical classifications of records in PubChem

- Includes links of different types

- Provides a way to find records with particular types of information

- Allows insightful comparisons to be made between sets of records

- Provides list import capability via Entrez (adding PubChem Search)

National Center for Biotechnology Information

## PubChem Classification Browser

Help

Browse PubChem data using a classification of interest, or search for PubChem records annotated with the desired classification/term (e.g., MeSH: phenylpropionates, or Gene Ontology: DNA repair). More...

Select classification
**MeSH**

Search selected classification by
**Keyword**    Enter desired search term    **Search**

Classification description (from MeSH)
MeSH (Medical Subject Headings) is the NLM controlled vocabulary thesaurus used for indexing articles for PubMed. More...

Data type counts to display
**None**  **Compound**  **Substance**  **PubMed**

Display zero count nodes?
**Yes**  **No**

Filter by Entrez History
#10 Search (#8 AND #9) (pubmed): 37435 results

### Browse MeSH Tree (filter applied ✕)

- ▼ MeSH Tree  ? ↗ 37,435
  - ▶ Analytical, Diagnostic and Therapeutic Techniques and Equipment Category ↗ 26,720
  - ▶ Anatomy Category ↗ 17,953
  - ▶ Anthropology, Education, Sociology and Social Phenomena Category ↗ 676
  - ▼ Chemicals and Drugs Category ↗ 37,435
    - ▶ Amino Acids, Peptides, and Proteins  ? ↗ 24,569
    - ▶ Biological Factors  ? ↗ 19,741
    - ▶ Biomedical and Dental Materials  ? ↗ 1,237
    - ▶ Carbohydrates  ? ↗ 17,361
    - ▶ Chemical Actions and Uses  ? ↗ 27,187
    - ▶ Complex Mixtures  ? ↗ 3,702
    - ▶ Enzymes and Coenzymes  ? ↗ 6,528
    - ▶ Heterocyclic Compounds  ? ↗ 8,321
    - ▶ Hormones, Hormone Substitutes, and Hormone Antagonists  ? ↗ 2,721
    - ▶ Inorganic Chemicals  ? ↗ 3,948
    - ▼ Lipids  ? ↗ 37,435

# Many Helpful Services and Functions

- Identifier Exchange Service
- Score Matrix Service
- Standardization Service
- BioActivity Dyad pages
- Entrez Indicies and Filters
- Bulk Download facilities
- AutoComplete Service
- PubChemRDF REST

## Bulk Download

PubChem data are available for bulk download on the PubChem FTP site
(ftp://ftp.ncbi.nlm.nih.gov/pubchem).

...a subset of PubChem records using the following services:

...load service (https://pubchem.ncbi.nlm.nih.gov/pc_fetch/pc_fetch.cgi)
...ne to download a list of compound or substance records in PubChem.  A list of
...d may be provided directly into the web page form or uploaded
...SID/CID per line with no heading). Alternatively, they may be provided by using
Entrez history, which stores a list of CIDs or SIDs returned from a previous Entrez search.  The records
can be exported in several formats, including SDF, PNG, SMILES, InChI, XML, and either text or binary
ASN.1.  The files may be optionally compressed in standard gzip (.gz) or bzip2 (.bz2) formats.

...Downloads through the structure download tool are limited to a maximum of 500,000 records per

Are specialized services in the cards?
To bulk download select sets?
To annotate?  To compare?
To subset and select?
To analyze with various methods?
Will users understand?

## Gene-CID dyad page

The Gene-CID dyad page shows the bioactivity data of a given compound record tested against a
particular gene target. For example, the following page presents the bioactivity data of CID 5328245
tested against GeneID 1956.

https://pubchem.ncbi.nlm.nih.gov/target/gene/1956#cid=5328245

This page can be accessed from the gene target page.
https://pubchem.ncbi.nlm.nih.gov/target/gene/EGFR/human#section=Tested-Compounds (click the
Structure or Activity column)

# Biologics

1.5M chemicals are 'biologics' in PubChem

- Glycan, amino acid, nucleic acid monomers

- Handles substitutions and chemical linkers

NCBI Glycans page

- Symbolic Nomenclature for Glycans (SNFG)

- Working with Glycans community

What is a glycan?

- WURCS collaboration

What is a biopolymer monomer?

- Pistoia Alliance HELM / EBI collaboration

## 2 Biologic Description



| SVG Image | H—Pro — His — Thr —Asn— Glu — Thr — Ser — Leu—OH (PO3H2) |
|---|---|
| IUPAC Condensed | H-Pro-His-Thr-Asn-Glu-Thr-Ser(PO3H2)-Leu-OH |
| Sequence | PHTNETXL |
| HELM | PEPTIDE1{P.H.T.N.E.T.[*C(=O)[C@H](COP(=O)(O)O)N* \|$_R2:;;;;;;;;;;;_R1$\|].L}$$$$ |
| IUPAC | L-prolyl-L-histidyl-L-threonyl-L-asparagyl-L-alpha-glutamyl-L-threonyl-O-phosphono-L-seryl-L-leucine |

▶ from PubChem

# New approaches needed with ultra-large databases

- Compact formats that resist enumeration
- Computational efficiency like never before
    - Computed property computation
    - How to establish links/relationships between entities?
- Rethinking the 'user' workflow
    - Machines/scripts vs. humans
- Cost efficiency
    - Expensive national resource or on each user desktop?

# Parting thoughts

- Very large databases (on the order of 1B-10B) are here

- Ultra-large databases (on the order of 100B-1TB) are just around the corner

- Much to do to scale

- Not every database needs or desires features found in PubChem but they will need some key features to be useful/relevant to end users

- Substantial investment in developing new algorithms, software, and rethinking databases will be necessary for practical utilization (features, cost, time)

Image credit:
https://www.myhubintranet.com/wp-content/uploads/2017/04/improving-employee-engagement.png

# PubChem Crew …

**Evan Bolton**

**Jie Chen**

**Tiejun Cheng**

**Asta Gindulyte**

**Jane He**

**Siqian He**

**Sunghwan Kim**

**Qingliang Li**

**Ben Shoemaker**

**Paul Thiessen**

**Bo Yu**

**Leonid Zaslavsky**

**Jian Zhang**

Special thanks to the NCBI Help Desk, especially Rana Morris, and past PubChem group members

# Special thanks

- Software collaborators
  - NextMove Software (Roger Sayle, John May) plus Daniel Lowe, Noel O'Boyle
  - Xemistry GmbH (Wolf D. Ihlenfeldt)
  - OpenEye Scientific Software
- PubChemRDF Collaborators
- BioHackathon (2014-2019) attendees and organizers
- All PubChem Contributors and Collaborators

-

[save the world •
love the people •
be happy]

Your
questions?